



INTRODUCTION

Artificial neural networks (ANNs) serve to process, learn, and predict information using layers of interconnected computational units resembling the neurons that comprise the biological nervous system. The ANN mimics the human brain in the ability of the latter to better adapt to experience, accommodate change, and learn from new perspectives, regardless of the unfamiliarity of the environment from which it learns. Although computers can perform many of the functions of the human brain, they lack the ability to learn from experience and to adapt to new perspectives, regardless of the unfamiliarity of the environment from which it learns. Although computers can perform many of the functions of the human brain, they lack the ability to learn from experience and to adapt to new perspectives, regardless of the unfamiliarity of the environment from which it learns.

Artificial neural networks (ANNs) serve to process, learn, and predict information using layers of interconnected computational units resembling the neurons that comprise the biological nervous system. The ANN mimics the human brain in the ability of the latter to better adapt to experience, accommodate change, and learn from new perspectives, regardless of the unfamiliarity of the environment from which it learns. Although computers can perform many of the functions of the human brain, they lack the ability to learn from experience and to adapt to new perspectives, regardless of the unfamiliarity of the environment from which it learns.

METHODOLOGY

Error Back-Propagation Training Algorithm (Levenberg-Marquardt)

- The network begins with arbitrarily weighted parameters
- A vector of x galaxy property input is fed into the network, which performs calculations to project an estimated value for the distance to the galaxy as the output
- The output vector is compared against the measured output distance; the mean squared difference between the two is noted as the error
- Using the gradient descent method, the network adjusts its parameter weights according to the direction of the steepest decreasing error
- After all galaxy property inputs are iterated through the network, the parameters will have adjusted in such a way that the network's outputted distance predictions become increasingly accurate of the correct distances and the performance error is minimized (see figure above)

Input Selection

- To determine where the network is making the mistakes, isolate each property input (i.e., u , s , r , i , z , petros0, petros90, or α and β) and run them separately through the network
- Input properties generating errors closest to the error produced by the inclusion of all nine inputs serve as the most significant contributors to the overall error
- The training algorithm is applied to all possible combinations of two property inputs, as opposed to individual inputs, in order to better approach the network's actual training potential,

Outlier Removal

- Eliminate outliers residing far from the rest of the data as well as the test data
- Idea: if the training data is reduced to a more general representation of the test set, the network will more easily recognize patterns in the data in order to extract a trend applicable to the test data
- Similar points reside in high density area
- Perpendicular bisectors (in red) designate voronoi cells
- points (in blue) designate voronoi cells
- points (in red) reside in low density area

- Voronoi cells define the space around data point closer to that point than to any other plotted point
- An outlier point far from the rest of the data set (in a low density area) would reside in a large cell
- If large voronoi cells correspond to poor network fittings, we can distinguish outliers by voronoi cell volume
- The network can then train to accommodate the likelihood of the higher error caused by an outlier

Improving Neural Network Generalization Ability Using Outlier Analysis & Voronoi Tessellation

M. Ho, D.M. McIntosh, A.N. Silvestri

NASA Summer High School Apprenticeship Research Program, Moffett Field, CA, USA
NASA Ames Research Center, Moffett Field, CA, USA

DATA

The data used in this study was obtained from the Sloan Digital Sky Survey (SDSS), which provides astronomers with what is currently the most extensive mapping of the universe, covering 25% of the sky and cataloging the spectral properties (e.g. luminosity, color, surface temperature) of over 100 million celestial objects.

Images generated by the SDSS are collected through five filters (ultraviolet, u , g , r , i , and z) that have respective wavelengths of 3540, 4790, 6222, 7632, and 9099 Å.

By measuring the photometric redshifts of the aforementioned wavelengths of a galaxy, astronomers can ascertain the extent to which the galaxy is receding from which the distance to the galaxy can be calculated.

Data collected for a small select group of galaxies (approximately 50,000) contains accurate measurements of the galaxies' redshifts in addition to measurements of their spectral properties.

The above data set containing both redshift measurements as well as spectral properties of the selected galaxies served as the training set for the purposes of this study; the data set containing only the spectral properties of a separate group of galaxies served as the test set.

FIGURES

Application of Training Algorithm on Altered Data Sets

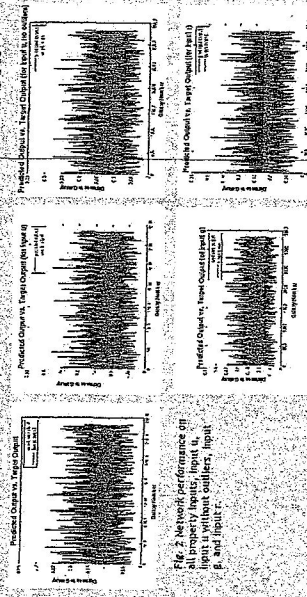


Fig. 2. Network performance on input B and input r.

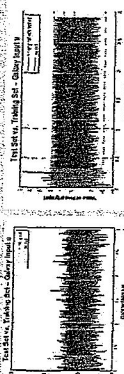


Fig. 3. Example of training set before and after outlier elimination. Data and the majority of the training data were removed; consequently, the outliers existing in the test set are not encompassed by the reduced training set.

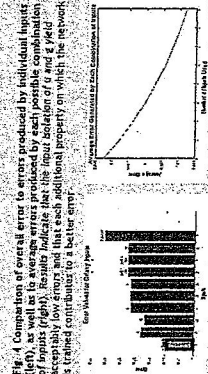


Fig. 4. Comparison of overall error to errors produced by individual inputs (left), as well as errors produced by each possible combination of inputs (right). Results indicate that the low error combination is trained contributors to a better error.

Voronoi Tessellation of Error Generated by Inputs u & r

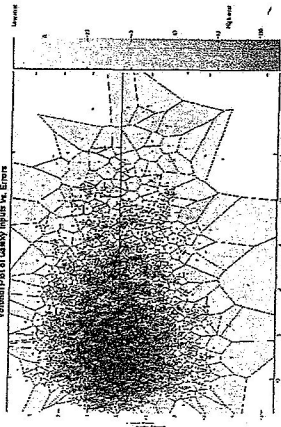


Fig. 1. The addition of the galaxies photographed above indicate that other objects within their vicinity.

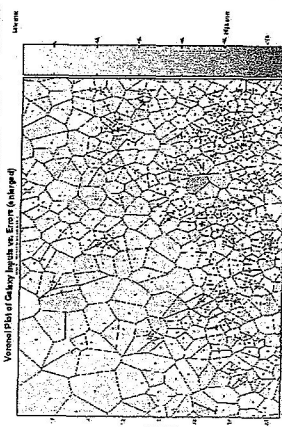


Fig. 5. A clear outlier residing in large voronoi cells exhibit surprising low errors, while several specific voronoi cells distributed throughout the training data exhibit particularly high errors and distinguish the areas of the training data in which the outliers appear to be randomly distributed throughout the tessellation.

CONCLUSIONS

Although the network's best performance resulted from the inclusion of all galaxy properties, the combination of galaxy property inputs u and r produced a notably low error that was considerably close to the overall error.

The network generated a lower error after training on data containing outliers, since the wide-ranging test set encompassed data anomalies to which the network's learned model could not have been applicable. By training the network on several outliers, the network grew more robust and improved its recognition of patterns appearing in an outlier training set that was difficult to generalize, thus reducing the risk of overfitting that would have arisen from training on an overly general or condensed data set.

Voronoi tessellation reveals that outliers with high voronoi cell volumes are associated with low to mid-range errors, which reinforces the notion that the network had already adjusted itself to place less voronoi volume on outliers, which may impact the training process for the network to further improve generalization ability.

For future research, the impact of including initially high or low errors as the network learns the relationship between distance error and voronoi cells of outliers should be studied, since the voronoi volumes corresponding to extreme errors which may impact the training process in the same way as outliers data points do.

References

1. Ho, M., D.M. McIntosh, A.N. Silvestri, "Improving Neural Network Generalization Ability Using Outlier Analysis & Voronoi Tessellation", NASA Summer High School Apprenticeship Research Program, Moffett Field, CA, USA, 2010.
2. Ho, M., D.M. McIntosh, A.N. Silvestri, "Improving Neural Network Generalization Ability Using Outlier Analysis & Voronoi Tessellation", NASA Summer High School Apprenticeship Research Program, Moffett Field, CA, USA, 2010.